

Decision Trees

Nipun Batra and teaching staff

July 16, 2025

IIT Gandhinagar

Discrete Input Discrete Output

The need for interpretability

Training Data

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Learning a Complicated Neural Network

Learnt Decision Tree



Medical Diagnosis using Decision Trees

Source: Improving medical decision trees by combining relevant health-care criteria

Optimal Decision Tree

Greedy Algorithm

Core idea: At each level, choose an attribute that gives **biggest estimated** performance gain!

Greedy!=Optimal

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insights: Problem is “easier” when there is lesser disagreement

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insights: Problem is “easier” when there is lesser disagreement
- Need some statistical measure of “disagreement”

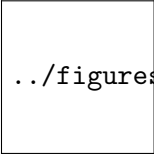
Statistical measure to characterize the (im)purity of examples

Entropy

Statistical measure to characterize the (im)purity of examples

$$H(X) = - \sum_{i=1}^n p(x_i) \log p(x_i)$$

Notebook: entropy.html



```
../figures/decision-trees/entropy.pdf
```

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Can we use Outlook as the root node?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Can we use Outlook as the root node?
- When Outlook is overcast, we always Play and thus no “disagreement”

Reduction in entropy by partitioning examples (S) on attribute A

$$\text{Gain}(S, A) \equiv \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$
 - Examples_v : subset of examples that $A = v$

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - $\text{Root} \leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$
 - Examples_v : subset of examples that $A = v$
 - If Examples_v is empty: add leaf with label = most common value of Target Attribute

ID3 (Examples, Target Attribute, Attributes)

- Create a root node for tree
- If all examples are $+/-$, return root with label = $+/-$
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $A \leftarrow$ attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $A = v$
 - Examples _{v} : subset of examples that $A = v$
 - If Examples _{v} is empty: add leaf with label = most common value of Target Attribute
 - Else: ID3 (Examples _{v} , Target attribute, Attributes - A)

Learnt Decision Tree

Root Node (empty)

Training Data

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Entropy calculated

We have 14 examples in S : 5 No, 9 Yes

$$\begin{aligned}\text{Entropy}(S) &= -p_{No} \log_2 p_{No} - p_{Yes} \log_2 p_{Yes} \\ &= -(5/14) \log_2(5/14) - (9/14) \log_2(9/14) = 0.94\end{aligned}$$

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Overcast	Yes
Rain	Yes
Rain	Yes
Rain	No
Overcast	Yes
Sunny	No
Sunny	Yes
Rain	Yes
Sunny	Yes
Overcast	Yes
Overcast	Yes
Rain	No

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3 No

Entropy =

$$-3/5\log_2(3/5) -$$

$$2/5\log_2(2/5) =$$

$$0.971$$

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3 No

$$\begin{aligned}\text{Entropy} &= \\ &= -3/5 \log_2(3/5) - \\ &= 2/5 \log_2(2/5) = \\ &= 0.971\end{aligned}$$

Outlook	Play
Overcast	Yes
Overcast	Yes
Overcast	Yes
Overcast	Yes

We have 4 Yes, 0 No

$$\text{Entropy} = 0$$

Information Gain for Outlook

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

We have 2 Yes, 3 No

$$\begin{aligned}\text{Entropy} &= \\ &= -3/5 \log_2(3/5) - \\ &= 2/5 \log_2(2/5) = \\ &= 0.971\end{aligned}$$

Outlook	Play
Overcast	Yes
Overcast	Yes
Overcast	Yes
Overcast	Yes

We have 4 Yes, 0 No

$$\text{Entropy} = 0$$

Outlook	Play
Rain	Yes
Rain	Yes
Rain	No
Rain	Yes
Rain	No

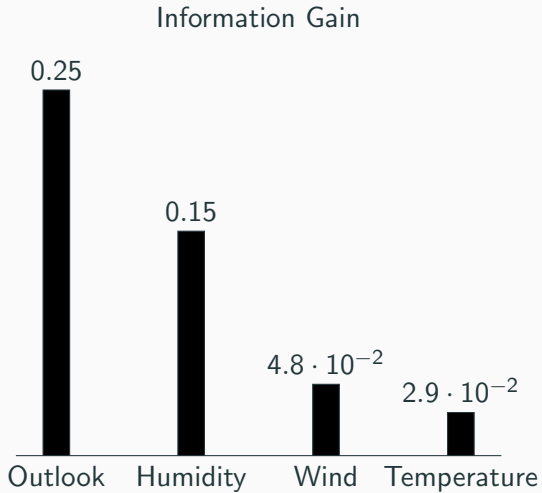
We have 3 Yes, 2 No

$$\begin{aligned}\text{Entropy} &= \\ &= -3/5 \log_2(3/5) - \\ &= 2/5 \log_2(2/5) = \\ &= 0.971\end{aligned}$$

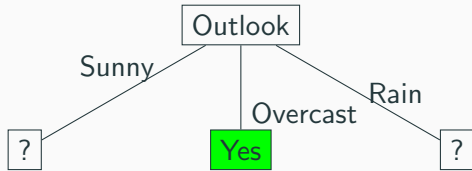
$$\text{Gain}(S, \text{Outlook}) = \text{Entropy}(S) - \sum_{v \in \{\text{Rain}, \text{Sunny}, \text{Windy}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\begin{aligned} \text{Gain}(S, \text{Outlook}) &= \text{Entropy}(S) - (5/14) * \text{Entropy}(S_{\text{Sunny}}) - \\ &\quad (4/14) * \text{Entropy}(S_{\text{Overcast}}) - (5/14) * \text{Entropy}(S_{\text{Rain}}) \\ &= 0.940 - 0.347 - 0.347 \\ &= 0.246 \end{aligned}$$

Information Gain



Learnt Decision Tree



Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- Gain($S_{\text{Outlook=Sunny}}$, Temp) = Entropy(3 Yes, 2 No) -
(2/5)*Entropy(2 No, 0 Yes) -(2/5)*Entropy(1 No, 1 Yes) -
(1/5)*Entropy(1 Yes)

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

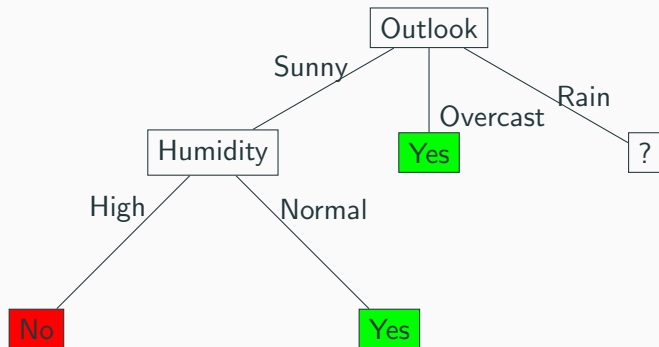
- $\text{Gain}(S_{\text{Outlook}=\text{Sunny}}, \text{Temp}) = \text{Entropy}(3 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(2 \text{ No}, 0 \text{ Yes}) - (2/5) * \text{Entropy}(1 \text{ No}, 1 \text{ Yes}) - (1/5) * \text{Entropy}(1 \text{ Yes})$
- $\text{Gain}(S_{\text{Outlook}=\text{Sunny}}, \text{Humidity}) = \text{Entropy}(3 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(2 \text{ Yes}) - (3/5) * \text{Entropy}(3 \text{ No}) \implies \text{maximum possible for the set}$

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- $\text{Gain}(S_{\text{Outlook}=\text{Sunny}}, \text{Temp}) = \text{Entropy}(3 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(2 \text{ No}, 0 \text{ Yes}) - (2/5) * \text{Entropy}(1 \text{ No}, 1 \text{ Yes}) - (1/5) * \text{Entropy}(1 \text{ Yes})$
- $\text{Gain}(S_{\text{Outlook}=\text{Sunny}}, \text{Humidity}) = \text{Entropy}(3 \text{ Yes}, 2 \text{ No}) - (2/5) * \text{Entropy}(2 \text{ Yes}) - (3/5) * \text{Entropy}(3 \text{ No}) \implies \textbf{maximum possible for the set}$
- $\text{Gain}(S_{\text{Outlook}=\text{Sunny}}, \text{Windy}) = \text{Entropy}(3 \text{ Yes}, 2 \text{ No}) - (3/5) * \text{Entropy}(2 \text{ No}, 1 \text{ Yes}) - (2/5) * \text{Entropy}(1 \text{ No}, 1 \text{ Yes})$

Learnt Decision Tree



Calling ID3 on (Outlook=Rain)

Day	Temp	Humidity	Windy	Play
D4	Mild	High	Weak	Yes
D5	Cool	Normal	Weak	Yes
D6	Cool	Normal	Strong	No
D10	Mild	Normal	Weak	Yes
D14	Mild	High	Strong	No

- The attribute Windy gives the highest information gain

Learnt Decision Tree



Prediction for Decision Tree

Just walk down the tree!



Prediction for Decision Tree

Just walk down the tree!



Prediction for <High Humidity, Strong Wind, Sunny Outlook, Hot Temp> is ?

Prediction for Decision Tree

Just walk down the tree!



Prediction for <High Humidity, Strong Wind, Sunny Outlook, Hot Temp> is ?
No

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit reached, the leaf node is assigned the “most” common occurring value in that path.

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit reached, the leaf node is assigned the “most” common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit reached, the leaf node is assigned the “most” common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Always predicting Yes

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit reached, the leaf node is assigned the “most” common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Always predicting Yes

What is depth-1 tree (no decision) for the examples?

Limiting Depth of Tree

Assuming if you were only allowed depth-1 trees, how would it look for the current dataset?

Apply the same rules, except when depth limit reached, the leaf node is assigned the “most” common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Always predicting Yes

What is depth-1 tree (no decision) for the examples?



Discrete Input, Real Output

Modified Dataset

Day	Outlook	Temp	Humidity	Wind	Minutes Played
D1	Sunny	Hot	High	Weak	20
D2	Sunny	Hot	High	Strong	24
D3	Overcast	Hot	High	Weak	40
D4	Rain	Mild	High	Weak	50
D5	Rain	Cool	Normal	Weak	60
D6	Rain	Cool	Normal	Strong	10
D7	Overcast	Cool	Normal	Strong	4
D8	Sunny	Mild	High	Weak	10
D9	Sunny	Cool	Normal	Weak	60
D10	Rain	Mild	Normal	Weak	40
D11	Sunny	Mild	High	Strong	45
D12	Overcast	Mild	High	Strong	40
D13	Overcast	Hot	Normal	Weak	35
D14	Rain	Mild	High	Strong	20

Measure of Impurity for Regression?

Measure of Impurity for Regression?

- Any guesses?

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $\text{MSE}(S) = 311.34$

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $MSE(S) = 311.34$
- Information Gain analogue?

Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- $MSE(S) = 311.34$
- Information Gain analogue?
- Reduction in MSE (weighted)

Gain by splitting on Wind

Wind	Minutes Played
Weak	20
Strong	24
Weak	40
Weak	50
Weak	60
Strong	10
Strong	4
Weak	10
Weak	60
Weak	40
Strong	45
Strong	40
Weak	35
Strong	20

$$\text{MSE}(S)=311.34$$

Wind	Minutes Played
Weak	20
Weak	40
Weak	50
Weak	60
Weak	10
Weak	60
Weak	40
Weak	35

Weighted

$$\text{MSE}(S_{\text{Wind}=\text{Weak}})=(8/14)*277=159)$$

Wind	Minutes Played
Strong	24
Strong	10
Strong	4
Strong	45
Strong	40
Strong	20

Weighted

$$\text{MSE}(S_{\text{Wind}=\text{Strong}})=(6/14)*218=93)$$

Notebook: [decision-tree-real-output.html](#)

```
../figures/decision-trees/discrete-input-real-output-lev
```


Real Input Discrete Output

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85
- Calculate the weighted impurity for each split

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85
- Calculate the weighted impurity for each split
- Choose the split with the lowest impurity

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 44
- LHS has 1 No and 0 Yes; RHS has 3 Yes and 2 No
- Entropy for LHS = 0, Entropy for RHS = 0.971
- Weighted Entropy = $0.971 * 5/6 = 0.808$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54
- LHS has 2 No and 0 Yes; RHS has 3 Yes and 1 No
- Entropy for LHS = 0, Entropy for RHS = 0.811
- Weighted Entropy = $0.811 * 4/6 = 0.541$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66
- LHS has 2 No and 1 Yes; RHS has 2 Yes and 1 No
- Entropy for LHS = 0.918, Entropy for RHS = 0.918
- Weighted Entropy = $0.918 \cdot 3/6 + 0.918 \cdot 3/6 = 0.918$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 76
- LHS has 2 No and 2 Yes; RHS has 1 Yes and 1 No
- Entropy for LHS = 1, Entropy for RHS = 1
- Weighted Entropy = $1 \cdot 4/6 + 1 \cdot 2/6 = 1$

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

Notebook: decision-tree-real-input-discrete-output.html

`../figures/decision-trees/real-ip-1.pdf`

Finding splits

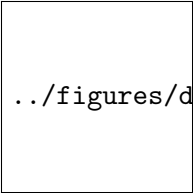
Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

Notebook: [decision-tree-real-input-discrete-output.html](#)

[../figures/decision-trees/real-ip-2.pdf](#)

Example (DT of depth 1)

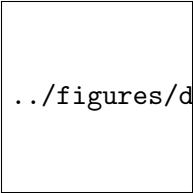
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-1.pdf
```

Example (DT of depth 2)

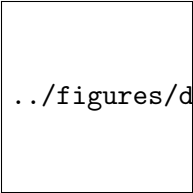
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-2.pdf
```

Example (DT of depth 3)

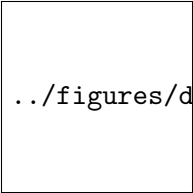
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-3.pdf
```

Example (DT of depth 4)

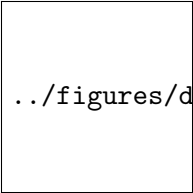
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-4.pdf
```


Example (DT of depth 5)

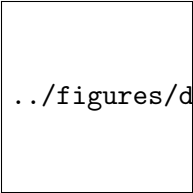
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-5.pdf
```

Example (DT of depth 6)

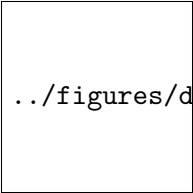
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-6.pdf
```

Example (DT of depth 7)

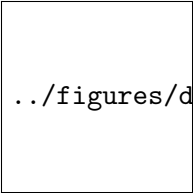
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-7.pdf
```

Example (DT of depth 8)

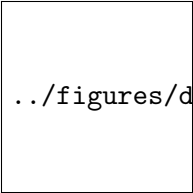
Notebook: [decision-tree-real-input-discrete-output.html](#)



```
../figures/decision-trees/dt-8.pdf
```

Example (DT of depth 9)

Notebook: [decision-tree-real-input-discrete-output.html](#)



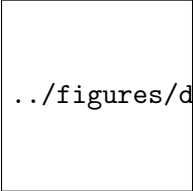
```
../figures/decision-trees/dt-9.pdf
```

Real Input Real Output

Example 1

Let us consider the dataset given below

Notebook: [decision-tree-real-input-real-output.html](#)

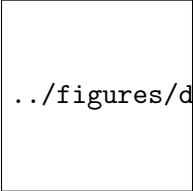


```
../figures/decision-trees/ri-ro-dataset.pdf
```

Example 1

What would be the prediction for decision tree with depth 0?

Notebook: [decision-tree-real-input-real-output.html](#)



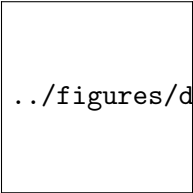
```
../figures/decision-trees/ri-ro-dataset.pdf
```


Example 1

Prediction for decision tree with depth 0.

Horizontal dashed line shows the predicted Y value. It is the average of Y values of all datapoints.

Notebook: [decision-tree-real-input-real-output.html](#)

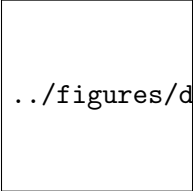


```
../figures/decision-trees/ri-ro-depth-0.pdf
```

Example 1

What would be the decision tree with depth 1?

Notebook: [decision-tree-real-input-real-output.html](#)

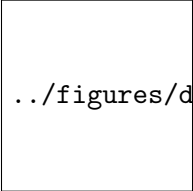


```
../figures/decision-trees/ri-ro-dataset.pdf
```

Example 1

Decision tree with depth 1

Notebook: [decision-tree-real-input-real-output.html](#)

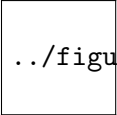


```
../figures/decision-trees/ri-ro-depth-1.pdf
```

Example 1

The Decision Boundary

Notebook: [decision-tree-real-input-real-output.html](#)

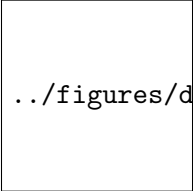


```
../figures/decision-trees/ri-ro-depth-1-sklearn.pdf
```

Example 1

What would be the decision tree with depth 2 ?

Notebook: [decision-tree-real-input-real-output.html](#)

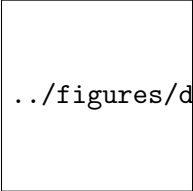


```
../figures/decision-trees/ri-ro-dataset.pdf
```

Example 1

Decision tree with depth 1

Notebook: [decision-tree-real-input-real-output.html](#)

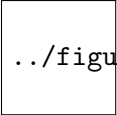


```
../figures/decision-trees/ri-ro-depth-2.pdf
```

Example 1

The Decision Boundary

Notebook: [decision-tree-real-input-real-output.html](#)



```
../figures/decision-trees/ri-ro-depth-2-sklearn.pdf
```

Objective function

Here, Feature is denoted by X and Label by Y .

Let the “decision boundary” or “split” be at $X = S$.

Let the region $X < S$, be region R_1 .

Let the region $X > S$, be region R_2 .

Objective function

Here, Feature is denoted by X and Label by Y .

Let the “decision boundary” or “split” be at $X = S$.

Let the region $X < S$, be region R_1 .

Let the region $X > S$, be region R_2 .

Then, let

$$C_1 = \text{Mean} (Y_i | X_i \in R_1)$$

$$C_2 = \text{Mean} (Y_i | X_i \in R_2)$$

Objective function

Here, Feature is denoted by X and Label by Y .

Let the “decision boundary” or “split” be at $X = S$.

Let the region $X < S$, be region R_1 .

Let the region $X > S$, be region R_2 .

Then, let

$$C_1 = \text{Mean } (Y_i | X_i \in R_1)$$

$$C_2 = \text{Mean } (Y_i | X_i \in R_2)$$

$$\text{Loss} = \sum_i ((Y_i - C_1 | X_i \in R_1)^2 + (Y_i - C_2 | X_i \in R_2)^2)$$

Objective function

Here, Feature is denoted by X and Label by Y .

Let the “decision boundary” or “split” be at $X = S$.

Let the region $X < S$, be region R_1 .

Let the region $X > S$, be region R_2 .

Then, let

$$C_1 = \text{Mean } (Y_i | X_i \in R_1)$$

$$C_2 = \text{Mean } (Y_i | X_i \in R_2)$$

$$\text{Loss} = \sum_i ((Y_i - C_1 | X_i \in R_1)^2 + (Y_i - C_2 | X_i \in R_2)^2)$$

Our objective is to minimize the loss and find

$$\min_S \sum_i ((Y_i - C_1 | X_i \in R_1)^2 + (Y_i - C_2 | X_i \in R_2)^2)$$

How to find optimal split “S”?

How to find optimal split “S”?

1. Sort all datapoints (X,Y) in increasing order of X .

How to find optimal split “S”?

1. Sort all datapoints (X,Y) in increasing order of X.
2. Evaluate the loss function for all

$$S = \frac{X_i + X_{i+1}}{2}$$

and then select the S with minimum loss.

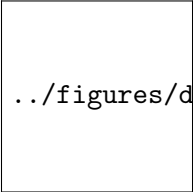
A Question!

Draw a regression tree for $Y = \sin(X)$, $0 \leq X \leq 2\pi$

A Question!

Dataset of $Y = \sin(X)$, $0 \leq X \leq 7$ with 10,000 points

Notebook: [decision-tree-real-input-real-output.html](#)

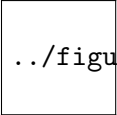


```
../figures/decision-trees/sine-dataset.pdf
```


A Question!

Regression tree of depth 1

Notebook: [decision-tree-real-input-real-output.html](#)

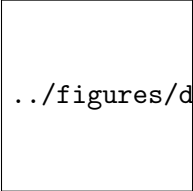


```
../figures/decision-trees/sine-depth-1-sklearn.pdf
```

A Question!

Decision Boundary

Notebook: [decision-tree-real-input-real-output.html](#)

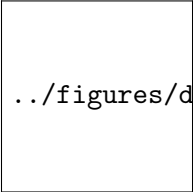


```
../figures/decision-trees/sine-depth-1.pdf
```

A Question!

Regression tree with no depth limit is too big to fit in a slide.
It has of depth 4. The decision boundaries are in figure below.

Notebook: [decision-tree-real-input-real-output.html](#)



```
../figures/decision-trees/sine-depth-4.pdf
```

Summary

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize “performance gain”
- Issues:
 - Can overfit easily!
 - Empirically not as powerful as other methods



`../figures/dt_weighted/fig1.pdf`



`../figures/dt_weighted/fig2.pdf`



../figures/dt_weighted/fig2.pdf

$$ENTROPY = -P(+)\cdot\log_2 P(+)-P(-)\cdot\log_2 P(-)$$

$$P(+)=\left(\frac{0.1+0.1+0.3}{1}\right)=0.5, P(-)=\left(\frac{0.3+0.1+0.1}{1}\right)=0.5$$

$$ENTROPY = E_s = -\frac{1}{2}\cdot\log_2\frac{1}{2}-\frac{1}{2}\cdot\log_2\frac{1}{2}=1$$

Weighted Entropy



`../figures/dt_weighted/fig3.pdf`

Candidate Line: $X1 = 4(X1^*)$



../figures/dt_weighted/fig4.pdf

Entropy of $X_1 \leq X_1^* = E_{S(X_1 < X_1^*)}$

$$P(+)=\left(\frac{0.1+0.1}{0.1+0.1+0.3}\right)=\frac{2}{5}$$

$$P(-)=\frac{3}{5}$$



../figures/dt_weighted/fig5.pdf

Entropy of $X_1 > X_1^* = E_{S(X_1 > X_1^*)}$

$$P(+)=\frac{3}{5}$$

$$P(-)=\frac{2}{5}$$



../figures/dt_weighted/fig5.pdf

$$IG(X_1 = X_1^*) = E_S - \frac{0.5}{1} \cdot E_{S(X_1 < X_1^*)} - \frac{0.5}{1} \cdot E_{S(X_1 > X_1^*)}$$